

000
001
002
003
004
005
006
007
008
009
010
011

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Vehicle Detection, Tracking, and Crash Prediction based on YOLO

Anonymous EECS 442 submission

Paper ID *****

Abstract

In the field of autonomous driving, detecting and predicting crashes may potentially help improve traffic situations and save lives. In this report, we explore a pipeline to track car trajectories at road crossing and predict potential accidents. Our project explores the YOLO model for object detection, algorithms for incorporating angle prediction for YOLO, vehicle identification, trajectory prediction, as well as crash risk assessments. We present an efficient pipeline to conduct crash detection using birds-eye-view input, contributing to the ongoing governmental efforts to improve road safety. Our project has the potential to be implemented in real-world road crossings to provide more agile emergency services in the event of accidents.

Keywords: YOLO, Object Detection and Tracking, Crash Prediction, Transfer Learning

1. Introduction

With modern autonomous driving technologies, we see a significant advancement in vehicular safety and traffic management, primarily through the use of bird's-eye view (BEV) cameras combined with machine learning algorithms. This paper focuses on enhancing the predictive capabilities of these technologies by adapting the YOLO (You Only Look Once) model to better interpret data provided by BEV cameras, thus improving its effectiveness in predicting vehicle trajectories and potential collision points.

Our research investigates modifications to the YOLO architecture to address the challenges of real-time, accurate path prediction and collision detection in dynamic traffic scenarios. Our study aims to contribute to safer autonomous driving systems and provide insights into the integration of advanced object detection technologies in traffic management.

2. Background

Safety measures, particularly in predicting and preventing potential traffic collisions become more and more significant these days. This research is motivated by the urgent need to reduce traffic accidents, which remain a leading cause of mortality globally. Utilizing computer vision to predict vehicular crashes before they occur can significantly contribute to the safety and efficiency of road traffic systems.

Bird's-eye view (BEV) cameras provide a comprehensive perspective of vehicular dynamics, offering a unique vantage point compared to traditional forward-facing vehicle cameras. This allows for more effective monitoring of vehicle trajectory, which are essential in predicting potential collision points. Our approach aims to explore these cameras' potential in the specific task of crash prediction.

The use of datasets like DOTA demonstrates the versatility of machine learning applications. However, adapting these datasets to train models for ground-level vehicular traffic analysis presents unique challenges, including the need for precise object detection and orientation. Our methodology addresses these challenges by refining detection algorithms to suit the specific vehicular situations and crash predictions.

In recent years, models such as YOLO have revolutionized object detection due to their speed and efficiency. However, the application of such models in dynamic, real-time environments like traffic systems requires significant modifications to accommodate the unique aspects of vehicular movements. This research expands on existing knowledge by adapting the YOLO architecture to detect oriented objects and predict their trajectories, thereby anticipating potential collisions more effectively. By enhancing the capabilities of existing computer vision techniques, we aim to contribute to the development of safer autonomous driving technologies and reduce the incidence of traffic accidents.

3. Methodology

To achieve the crash prediction feature we planned, there are a few key steps to conduct: 1. car detection 2. car identi-

108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
fication 3. trajectory prediction 4. determining crashes. Below we discuss implementation details and considerations for each part.

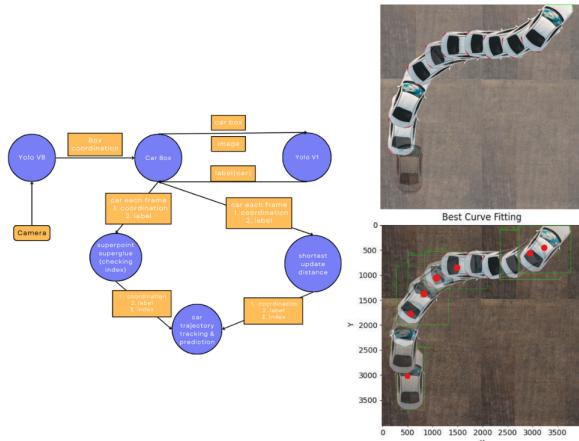


Figure 1. ROS Package Pipeline

Based on the ROS package we created, the user can input the RGB images from the camera, and the ROS package will finally output the car position, the car trajectory and the future driving trajectory as shown in Fig. 1.

3.1. Dataset

Our primary goal is to predict car crashes with a birds-eye view camera live feed. Naturally, the dataset we incorporate in our model training and system design should primarily include birds-eye-view (BEV) images. the DOTA [10] dataset is used for this purpose.

One of the datasets we tried to inspect is the DOTA dataset v1.0, a dataset dedicated to detecting objects from satellite and aerial images. The 1.0 version gives 20 different object categories including genres like basketball court or bridge, but for tracking purposes, only five of them are objects that can move: large vehicle, small vehicle, helicopter, plane, and ship should be taken out and inspected. Thus, one of the good ways to redirect the dataset to our purposes is to only filter out the tags needed. Besides, we will have to remove the images from the dataset that does not contain any labels that we would want.

Another thing we must take into account is the fact that DOTA dataset is annotated under the format “ $x_1, y_1, x_2, y_2, x_3, y_3, x_4, y_4$, category, difficult”, which are the four points of the quadrilateral for the label. For the purpose of OBB detection, we might have to translate the four points into $xywh\theta$ formatting during data loading. If we are not interested in orientation in this case, we will have to translate the coordinates into a two-point label system by taking the min and max of x and y coordinates respectively to align with the output from YOLO.

Furthermore, one additional limitation of the dataset is that the size and shape of the input image can vary by a great extent and it would be fundamentally inefficient to shrink the entire image into a 448x448 image as a car might become less than 10 pixels as shown in Fig. 2.

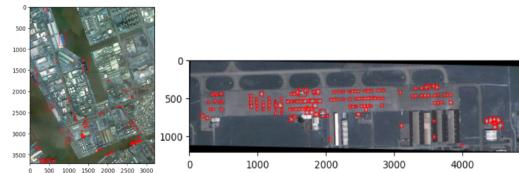


Figure 2. Example Images from the DOTA Dataset

Instead, though it slightly defeats the definition of the YOLO, we can essentially crop the images to smaller samples of it so that smaller objects are easier to be detected given limited input dimensions. Even given the upsampling structures as seen in YOLO v3 [6], the objects shown in the network will still have too few pixels to be detectable even by human eyes. Else, if the width and height of the images are too different, the objects can be distorted as a result.

To fix this, we can crop the images into smaller rectangular sizes like (448, 448) while retaining all the labels, filter out the images we need, and the model will have an easier time recognizing the patterns given the fixed size model. Here are some examples of the cropped samples from the original training set:

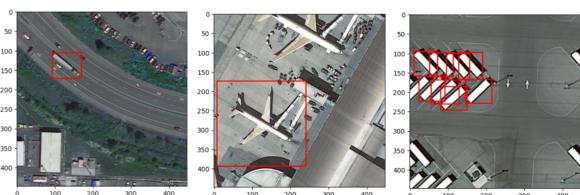


Figure 3. Examples of the Cropped Samples from DOTA

This is an example of how to adapt datasets into different purposes and some considerations before putting everything into training the model. However, during our initial attempts to actually train the network, images from DOTA is still too big for our purposes. Thus, we implemented transfer learning using PASCAL_VOC [2] dataset.

3.2. Car Detection

We chose the YOLO v1 [5] model as the base algorithm for car detection. Due to our more narrow use case, we modified the model output layer to predict fewer classes (5, rather than 20). We also adopted transfer learning techniques, hoping to generalize the model to be resilient in the case of extraordinary input. Finally, due to the oriented nature of the DOTA dataset, we modified the model output

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

216 layer definition as well as the loss function to correctly calculate
 217 the orientation loss.
 218

219 3.2.1 Detection Model Architecture 220

221 Table 1. Neural Network Architecture
 222

Type	Filters	Size/Stride	Output
Conv	64	7 x 7 / 2	224 x 224
Max Pool		2 x 2 / 2	112 x 112
Conv	192	3 x 3 / 1	112 x 112
Max Pool		2 x 2 / 2	56 x 56
		1 x	
Conv	128	1 x 1 / 1	56 x 56
Conv	256	3 x 3 / 1	56 x 56
Conv	256	1 x 1 / 1	56 x 56
Conv	512	3 x 3 / 1	56 x 56
Max Pool		2 x 2 / 2	28 x 28
		4 x	
Conv	256	1 x 1 / 1	28 x 28
Conv	512	3 x 3 / 1	28 x 28
Conv	512	1 x 1 / 1	28 x 28
Conv	1024	3 x 3 / 1	28 x 28
Max Pool		2 x 2 / 2	14 x 14
		2 x	
Conv	512	1 x 1 / 1	14 x 14
Conv	1024	3 x 3 / 1	14 x 14
Conv	1024	3 x 3 / 1	14 x 14
Conv	1024	3 x 3 / 2	7 x 7
Conv	1024	3 x 3 / 1	7 x 7
Conv	1024	3 x 3 / 1	7 x 7
FC		4096	4096
Dropout 0.5			4096
FC		7 x 7 x 30	7 x 7 x 30
FC		496	496
Dropout 0.5			496
FC		7 x 7 x 15	7 x 7 x 15
LeakReLU		7 x 7 x 15	7 x 7 x 15

260 The original YOLO v1 architecture we implemented can
 261 detect some labels such as person, cat, dog kind of good.
 262 However, for the car, mbike, and bus, it has a really bad
 263 detection, so in order to solve this problem, we will use the
 264 transfer learning method. Tab. 1.
 265

266 3.2.2 Transfer Learning 267

268 We use the YOLO v1 Darknet architecture with 24 layers
 269 combined with transfer learning to train our custom com-

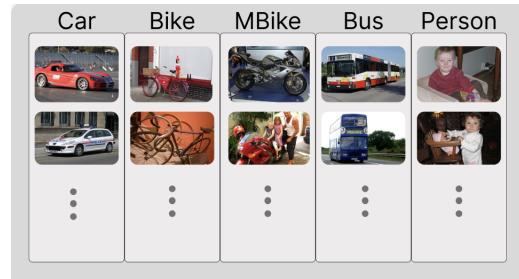
tuned parameters	values
pretrained_dataset	heavy_3000_20.csv
batch size	64
Learning Rate	$f(x) = \begin{cases} 2e^{-5} & \text{if mAP} < 0.65 \\ 2e^{-6} & \text{if mAP} \geq 0 \end{cases}$
transfer learning learnable/freeze	<ol style="list-style-type: none"> 1. change the last 2 fully connection layers to fit in the 5 labels 2. freeze the first 18 layers. 3. make the last 6 CNN and 2 new fully connection layers be learnable.
transfer learning dataset	transfer_1000_5.csv

291 Table 2. Transfer Learning Parameters
 292

293 puter vision model, specifically for detecting **cars**, **buses**,
bikes, **mbike**, and **persons**. Fig. 4The detailed architecture
 294 is specified in Tab. 1.

295 We initially utilized the original YOLO v1 architecture
 296 along with 3000 randomly selected images from 20 labels
 297 to train the model on drawing bounding boxes and identifying
 298 large objects within each frame. However, given our
 299 project's emphasis on predicting car trajectories and issuing
 300 collision warnings, we use transfer learning to let the model
 301 specifically learn on these five labels. The parameters tun-
 302 ning and transfer learning details are in Tab. 2

303 To accommodate the prediction of five classes, we mod-
 304 ify the output layers from 7x7x30 to 7x7x15 as 15 =
 305 $2 * 5 * |\text{bbox_params}| + 5 * |\text{labels}|$.
 306



311 Figure 4. Five Classes
 312

324

3.2.3 Complex (Oriented) YOLO Loss

To detect car bounding boxes that are actually oriented, we incorporate oriented bounding box notation, where a bounding box is described not by (x, y, w, h) in YOLO [5], but (x, y, w, h, θ) . Essentially, besides the center, width, height of a bounding box, an additional parameter θ is added to describe the heading of the bounding box. This modification to the YOLO loss is presented in Complex YOLO [9]. Complex-YOLO proposed L_{Euler} , the loss related to the bounding box angle, described with complex numbers. Formally, this part of the loss is defined as

$$\begin{aligned}\mathcal{L}_{Euler} &= \lambda \sum_{S^2} \sum_B \mathbb{1}^{obj} |e_{pred} - e_{gt}|^2 \\ &= \lambda \sum_{S^2} \sum_B \mathbb{1}^{obj} [(r_{pred} - r_{gt})^2 + (i_{pred} - i_{gt})^2]\end{aligned}$$

where e_{pred}, e_{gt} are complex numbers on the unit circle representing the angle of predicted and ground-truth bounding box, respectively, and r, i are their real and imaginary parts, respectively. Thus, the Complex YOLO loss is defined as

$$\mathcal{L}_{Comp} = \mathcal{L}_{Yolo} + \mathcal{L}_{Euler}$$

We ultimately ran out of time to modify our model output and train on this modified loss, but it is certainly worth exploring as an extension to the traditional YOLO algorithm.

3.2.4 Oriented Bounding Box Detection and IOU

One part of the YOLO loss requires the computation of the Intersection-Over-Union (IOU) between the predicted bounding box and the ground truth. The IOU between two bounding boxes is defined as

$$IOU = \frac{Area(B_{pred} \cap B_{gt})}{Area(B_{pred} \cup B_{gt})}$$

While it is straightforward for axis-aligned boxes, computing the IOU of two oriented bounding boxes are much more complex. We developed **PolygonIntersect()** to accomplish this objective. This algorithm computes the polygon intersection of two arbitrary convex quadrilateral in $O(n^2)$ time.

While traditional computer graphics often use the ray casting algorithm [8] to check whether a point is within a polygon, we observed that for a point to be inside a convex polygon, the vector it forms with the vertices of the polygon must have angles of the same sign with respect to corresponding edges. This information is easily attainable through a 2D cross product (or determinant). This is the theoretical foundation of Algorithm 2. While the complexity of Algorithm 2 remains $O(n)$, same as the ray-casting algorithm, this could be more computationally efficient when implemented in vectorized form in a computer system.

Algorithm 1 Polygon Intersection

```

1: function INTERSECT( $G_1 = (V_1, E_1), G_2 = (V_2, E_2)$ )
2:   Initialize  $I \leftarrow []$ ,  $c \leftarrow Arr[|E_1|, |E_2|]$ 
3:   for  $(e_1, e_2) \in (E_1, E_2)$  do
4:     if  $e_1, e_2$  not parallel or overlapping then
5:        $c_{ij} \leftarrow \text{Intersection}(e_1, e_2)$ 
6:     end if
7:   end for
8:    $I \leftarrow \{c_{ij} : c \text{ lies on } e_i, e_j\}$ 
9:    $P_{1in2} = \{p_1 \in V_1 : \text{Within}(G_2, p_1)\}$ 
10:   $P_{2in1} = \{p_2 \in V_2 : \text{Within}(G_1, p_2)\}$ 
11:   $I \leftarrow I \cup P_{1in2} \cup P_{2in1}$ 
12: return  $I$ 
13: end function
```

Algorithm 2 Within

```

1: function WITHIN( $G = (V, E), p$ )
2:    $S \leftarrow Arr[|E|]$ 
3:   for  $(e \in E \text{ do}$ 
4:      $x_1, y_1 \leftarrow \text{begin}(e), x_2, y_2 \leftarrow \text{end}(e)$ 
5:      $x_p, y_p \leftarrow p$ 
6:      $vec_e \leftarrow (x_2 - x_1, y_2 - y_1)$ 
7:      $vec_p \leftarrow (x_p - x_1, y_p - y_1)$ 
8:      $S_i \leftarrow \text{sign}(vec_p \times vec_e)$ 
9:   end for
10:  return ( $S_i >= 0, \forall i \in |E| \vee (S_i <= 0, \forall i \in |E|)$ 
11: end function
```

Upon finding the description of the polygon, we then proceed to finding the area of the polygon intersection. We implement the Area of Polygon algorithm [4] to compute the area of an arbitrary (non-self intersecting) polygon in $O(|V|)$ time.

This suite of algorithms allows us to compute the areas, thus the IOU, of arbitrary convex quadrilateral intersects, which is perfect for the task and dataset we chose. Results and visualization of this algorithm is presented in the Results section.

Unfortunately, we ultimately were not able to successfully modify our model and corresponding loss function to train the model properly. Thus, we only presented results from the transfer learning model using the PASCAL_VOC dataset and the original YOLO loss.

3.2.5 Detection Classification

Due to the limitations of YOLO v1 that we discovered towards the end of the project, we were unable to use YOLO v1 to detect small objects from a large BEV image. To accommodate this, we tried to incorporate the YOLO v8 [3] model into the detection pipeline. While the pretrained YOLO v8 model fails to classify top-down views of cars

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432 most of the time, our transfer-learning model is able to successfully detect this if the image is cropped out.
 433
 434

435 Thus, for our final version of the prototype, we first
 436 detect objects with the YOLO v8 model, which most of
 437 the time correctly bounds the box (but produces a wrong
 438 label, like phone), crops this part of the image out, and
 439 then feed this into our YOLO v1 transferred model to de-
 440 termine whether the cropped image is a car or not. This
 441 method bears much improvement but it demonstrates that
 442 our transfer-learning demonstrates an improvement over
 443 even the vanilla YOLO v8 on specific tasks like top-down-
 444 view vehicle detection.
 445

3.3. Trajectory Prediction

446 To predict the trajectory of all cars detected, we must
 447 conduct two tasks: 1. Identify each car and make sure the
 448 waypoint is correctly assigned to the same car if previously
 449 detected. 2. Predict the trajectory given previous waypoints
 450 detected.
 451

3.3.1 Car Identification

452 We combined k-Nearest Neighbor with the SuperPoint
 453 [1]-SuperGlue [7] feature-fitting algorithm to select the car
 454 that is not only close to our predicted point, but also resem-
 455 bles the original car the most in the new image, as seen in
 456 Algorithm 3.
 457

Algorithm 3 CarID

```
1: function CARID( $C \in \mathbb{R}^{n \times 2}, P \in \mathbb{R}^{n \times 2}, I$ )
2:    $ID \leftarrow Arr[n]$ 
3:   for  $c_i \in C$  do
4:     Candidates =  $kNN(C_i, P)$ 
5:      $ID[i] \leftarrow \text{argmax}_i SPSG(\text{Candidates}[i], I)$ 
6:   end for
7:   return ID
8: end function
```

3.3.2 Trajectory Prediction

472 We attempted to fit polynomial lines to known way-
 473 points of a car in order to predict its trajectory in the near-
 474 future. To achieve this, we first find the best fit polynomial
 475 with order no more than 3 using the mean squared loss:
 476

$$\text{argmin}_i \frac{1}{n} \sum_{i=1}^n \frac{(y_i - \text{BestFit}(P, i)(x_i))^2}{2}$$

482 Then, the parametrized best fit polynomial is used to
 483 predict the next waypoint based on the car's current veloci-
 484 ty. We use NumPy to compute our best-fit polynomials on
 485 the past 4 waypoints (fewer if car just entered frame).

3.4. Crash Detection

486 To compute whether two cars may crash, we uses an al-
 487 gorithm structure as shown in Algorithm 4 to detect whether
 488 the distance between the cars, based on the predicted path,
 489 will be lower than a danger threshold value. Way-points are
 490 sampled from the provided polynomials at fixed time inter-
 491 vals for the near future, and distance values are computed
 492 to determine the risk of crashing.
 493

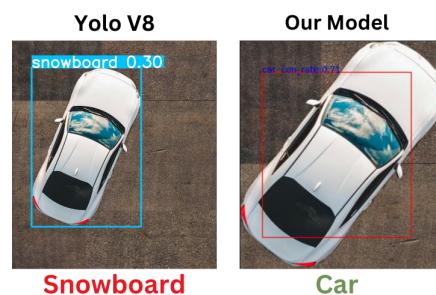
Algorithm 4 Crash Detection

```
1: function CRASHDETECT( $Poly_1, Poly_2$ )
2:    $P_1, P_2 \leftarrow \text{Sample}(Poly_1), \text{Sample}(Poly_2)$ 
3:    $D = \min_t |P_{1,t} - P_{2,t}|^2$ 
4:   return  $D \leq \text{DangerThreshold}$ 
5: end function
```

4. Results

4.1. Car Detection Model

505 YOLO v1 cannot detect small objects in the frame, at
 506 the same time, the YOLO v8 model cannot recognize the
 507 top-down view of the car as a car. So we finally use YOLO
 508 v8 to generate the box, and then based on the position of
 509 the bounding box, extract the cars from the image and pass
 510 them into the YOLO v1 model for recognition, ultimately
 511 determining their labels.Fig. 5



512 Figure 5. YOLO V8 VS Our Result

4.1.1 Model Training and Transfer Learning

528 After transfer learning, our model can effectively rec-
 529 ognize the five desired labels, while ignoring any objects
 530 outside of these five categories. The comparison of per-
 531 formance before and after transfer learning can be visually
 532 observed in Fig. 6.

533 In Fig. 6, the model with pre-trained weights as shown
 534 on the lower portion of the figure produced observably
 535 more accurate results in recall, classification accuracy, and
 536 bounding box accuracy. Given limited hardware capacities,
 537 transfer learning will be a much easier approach to

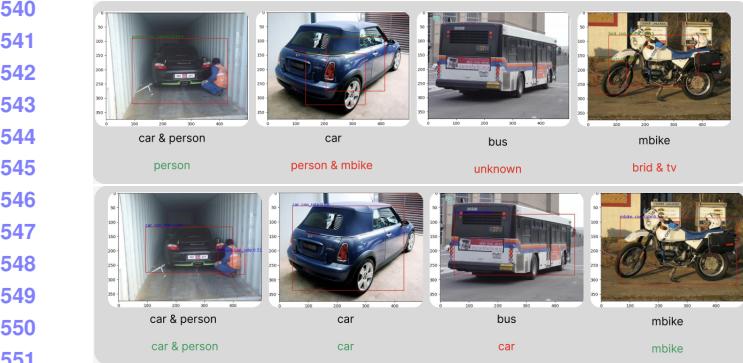


Figure 6. Transfer Learning Results Comparison

4.1.2 Oriented Bounding Box Detection

For Oriented Bounding Box (OBB), our oriented_iou implementation successfully solves for intersects between bounding boxes, their respective areas, as well as the IOU values. It also successfully identifies cases of inclusive, disjoint, and axis-aligned boxes, as shown in Fig. 7.

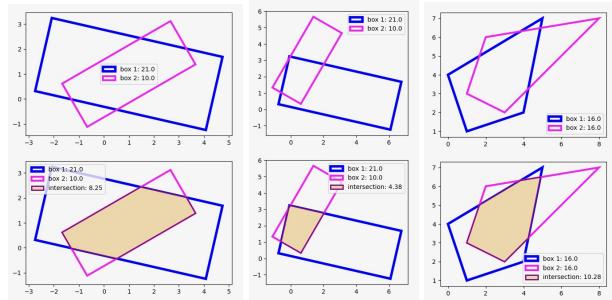


Figure 7. Oriented Bounding Box Detection and Calculation (left: two oriented quadrilaterals, right: detected intersection polygon)

4.2. Car Identity

As shown in Fig. 8, by comparing two photos taken at different times of the same intersection, we employ the SPSG algorithm to perform feature point matching, using white cars and gray cars as examples. We can observe that the white car in the bottom right corner of the second photo belongs to the same vehicle as the white car in the bottom right corner of the first photo, rather than being the same as the gray car in the top left corner of the first photo.

However, our model may also have vulnerabilities in cases where two cars of the same make and color, traveling at nearly identical speeds, intersecting trajectories, and both the KNN and SPSG algorithms may fail simultaneously, leading to errors in car index recognition. However, such instances did not occur in our subsequent test cases. The probability of such occurrences is extremely low.



Figure 8. SPSG matching result

4.3. Trajectory Prediction

Our algorithm successfully demonstrated resilience against the extraordinary scenarios discussed in the methods section. In most cases, KNN can correctly identify the car index, but in some extreme situations, it may make the wrong detection and prediction. For example, in Fig. 9, when it comes to 4th frame, relying solely on the KNN algorithm would incorrectly identify the position of the car in the top left corner as that of the car in the bottom right corner. After introducing the SPSG algorithm, discretion between rooftop patterns of two cars result in correct identification, as seen in Fig. 10.

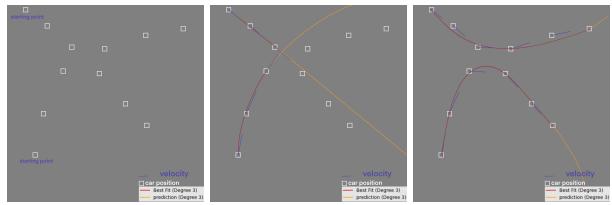


Figure 9. Wrong Prediction

Fig. 11 is an example of six frames where two cars are simultaneously traveling in different directions. Through our detection, recognition, and trajectory prediction pipeline, we can successfully detect the positions of vehicles and predict their trajectories in the near future.

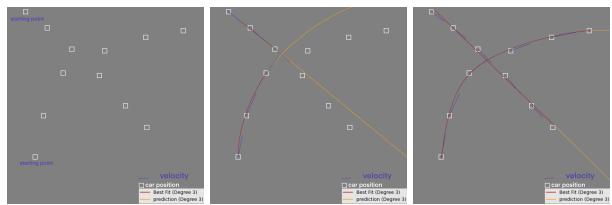


Figure 10. Correct Prediction

5. Conclusion

This paper presents attempts to explore in the task of autonomous driving, focusing on vehicle detection, tracking, and crash prediction using an adapted version of YOLO

594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

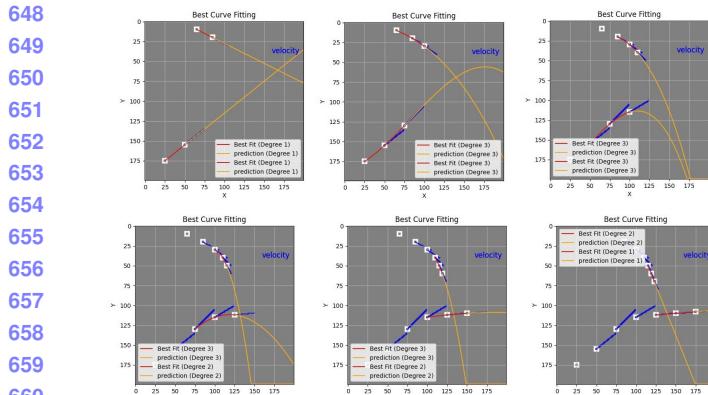


Figure 11. Trajectory Prediction Plots

model. Our modifications to the YOLO architecture, particularly the inclusion of oriented bounding boxes and the integration of a complex YOLO loss, have demonstrated the potential to enhance the accuracy of vehicle trajectory prediction and crash detection.

Moreover, transfer learning has refined our model's performance without relying on extensive datasets, allowing it to more accurately identify and track vehicles under road conditions. The implementation of algorithms for oriented bounding box detection and the calculation of intersection-over-union for complex polygons have provided robust tools for enhancing detection accuracy. Future work will focus on further optimizing the detection speed and reliability, as well as expanding the model's capabilities to include more diverse traffic scenarios and environmental conditions. Please refer to https://github.com/yuzhench/EECS442_final_project.git for more information.

References

- [1] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superpoint: Self-supervised interest point detection and description. *CoRR*, abs/1712.07629, 2017. 5
- [2] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision*, 88(2):303–338, June 2010. 2
- [3] Glenn Jocher, Ayush Chaurasia, and Jing Qiu. Ultralytics YOLO, Jan. 2023. 4
- [4] John Page. Algorithm to find the area of a polygon. <https://www.mathopenref.com/coordpolygonarea2.html>. Accessed: 2024-4-29. 4
- [5] Joseph Redmon, Santosh Kumar Divvala, Ross B. Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. *CoRR*, abs/1506.02640, 2015. 2, 4
- [6] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *CoRR*, abs/1804.02767, 2018. 2
- [7] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superglue: Learning feature matching with graph neural networks. *CoRR*, abs/1911.11763, 2019. 5
- [8] M. Shimrat. Algorithm 112: Position of point relative to polygon. *Commun. ACM*, 5(8):434, aug 1962. 4
- [9] Martin Simon, Stefan Milz, Karl Amende, and Horst-Michael Gross. Complex-yolo: Real-time 3d object detection on point clouds. *CoRR*, abs/1803.06199, 2018. 4
- [10] Gui-Song Xia, Xiang Bai, Jian Ding, Zhen Zhu, Serge J. Beßongie, Jiebo Luo, Mihai Datcu, Marcello Pelillo, and Liangpei Zhang. DOTA: A large-scale dataset for object detection in aerial images. *CoRR*, abs/1711.10398, 2017. 2